

CLAIMS

1. A method comprising:
receiving a frame of content;
automatically detecting a candidate area for a new face region in the frame;
using one or more hierarchical verification levels to verify whether a human face is in the candidate area;
indicating that the candidate area includes a face if the one or more hierarchical verification levels verify that a human face is in the candidate area;
and
using a plurality of cues to track each verified face in the content from frame to frame.
2. A method as recited in claim 1, wherein the frame of content comprises a frame of video content.
3. A method as recited in claim 1, wherein the frame of content comprises a frame of audio content.
4. A method as recited in claim 1, wherein the frame of content comprises a frame of both video and audio content.
5. A method as recited in claim 1, further comprising repeating the automatic detecting in the event tracking of a verified face is lost.

1 6. A method as recited in claim 1, wherein receiving the frame of
2 content comprises receiving a frame of video content from a video capture device
3 local to a system implementing the method.

4
5 7. A method as recited in claim 1, wherein receiving the frame of
6 content comprises receiving the frame of content from a computer readable
7 medium accessible to a system implementing the method.

8
9 8. A method as recited in claim 1, wherein detecting the candidate area
10 for the new face region in the frame comprises:

11 detecting whether there is motion in the frame and, if there is motion in the
12 frame, then performing motion-based initialization to identify one or more
13 candidate areas;

14 detecting whether there is audio in the frame, and if there is audio in the
15 frame, then performing audio-based initialization to identify one or more
16 candidate areas; and

17 using, if there is neither motion nor audio in the frame, a fast face detector
18 to identify one or more candidate areas.

19
20 9. A method as recited in claim 1, wherein detecting the candidate area
21 for the new face region in the frame comprises:

22 determining whether there is motion at a plurality of pixels on a plurality of
23 lines across the frame;

24 generating a sum of frame differences for each possible segment of each of
25 the plurality of lines;

1 selecting, for each of the plurality of lines, the segment having the largest
2 sum;
3 identifying a smoothest region of the selected segments;
4 checking whether the smoothest region resembles a human upper body; and
5 extracting, as the candidate area, the portion of the smoothest region that
6 resembles a human head.

7
8 **10.** A method as recited in claim 9, wherein determining whether there
9 is motion comprises:

10 determining, for each of the plurality of pixels, whether a difference
11 between an intensity value of the pixel in the frame and an intensity value of a
12 corresponding pixel in one or more other frames exceeds a threshold value.
13

14 **11.** A method as recited in claim 1, wherein the one or more hierarchical
15 verification levels include a coarse level and a fine level, wherein the coarse level
16 can verify whether the human face is in the candidate area faster but with less
17 accuracy than the fine level.
18

19 **12.** A method as recited in claim 1, wherein using one or more
20 hierarchical verification levels comprises, as one of the levels of verification:

21 generating a color histogram of the candidate area;
22 generating an estimated color histogram of the candidate area based on
23 previous frames;
24 determining a similarity value between the color histogram and the
25 estimated color histogram; and

1 verifying that the candidate area includes a face if the similarity value is
2 greater than a threshold value.

3
4 **13.** A method as recited in claim 1, wherein indicating that the candidate
5 area includes a face comprises recording the candidate area in a tracking list.

6
7 **14.** A method as recited in claim 13, wherein recording the candidate
8 area in the tracking list comprises accessing a record corresponding to the
9 candidate area and resetting a time since last verification of the candidate.

10
11 **15.** A method as recited in claim 1, wherein the one or more hierarchical
12 verification levels include a first level and a second level, and wherein using the
13 one or more hierarchical verification levels to verify whether the human face is in
14 the candidate area comprises:

15 checking whether, using the first level verification, the human face is
16 verified as in the candidate area; and

17 using the second level verification only if the checking indicates that the
18 human face is not verified as in the candidate area by the first level verification.

19
20 **16.** A method as recited in claim 1, wherein using one or more
21 hierarchical verification levels comprises:

22 using a first verification process to determine whether the human head is in
23 the candidate area; and

1 if the first verification process verifies that the human head is in the
2 candidate area, then indicating the area includes a face, and otherwise using a
3 second verification process to determine whether the human head is in the area.
4

5 17. A method as recited in claim 16, wherein the first verification
6 process is faster but less accurate than the second verification process.
7

8 18. A method as recited in claim 1, wherein the plurality of cues include
9 foreground color, background color, edge intensity, motion, and audio.
10

11 19. A method as recited in claim 1, wherein using the plurality of cues
12 to track each verified face comprises, for each face:

13 predicting where a contour of the face will be;
14 encoding a smoothness constraint that penalizes roughness;
15 applying the smoothness constraint to a plurality of possible contour
16 locations; and

17 selecting the contour location having the smoothest contour as the location
18 of the face in the frame.
19

20 20. A method as recited in claim 19, wherein the smoothness constraint
21 includes contour smoothness.
22

23 21. A method as recited in claim 19, wherein the smoothness constraint
24 includes both contour smoothness and region smoothness.
25

1 **22.** A method as recited in claim 19, wherein encoding the smoothness
2 constraint comprises generating Hidden Markov Model (HMM) state transition
3 probabilities.

4
5 **23.** A method as recited in claim 19, wherein encoding the smoothness
6 constraint comprises generating Joint Probability Data Association Filter (JPDAF)
7 state transition probabilities.

8
9 **24.** A method as recited in claim 19, wherein using the plurality of cues
10 to track each verified face further comprises, for each face:

11 adapting the predicting for the face in subsequent frames to account for
12 changing color distributions.

13
14 **25.** A method as recited in claim 19, wherein using the plurality of cues
15 to track each verified face further comprises, for each face:

16 adapting the predicting for the face in subsequent frames based on one or
17 more cues observed in the frame.

18
19 **26.** A method as recited in claim 1, wherein using the plurality of cues
20 to track each verified face comprises, for each face:

21 accessing a set of one or more feature points of the face;

22 analyzing the frame to identify an area that includes the set of one or more
23 feature points;

24 encoding a smoothness constraint that penalizes roughness;
25

1 applying the smoothness constraint to a plurality of possible contour
2 locations; and

3 selecting the contour location having the smoothest contour as the location
4 of the face in the frame.

5
6 **27.** A method as recited in claim 1, wherein using the plurality of cues
7 to track each verified face comprises concurrently tracking multiple possible
8 locations for the face from frame to frame.

9
10 **28.** A method as recited in claim 27, further comprising using a
11 multiple-hypothesis tracking technique to concurrently track the multiple possible
12 locations.

13
14 **29.** A method as recited in claim 27, further comprising using a particle
15 filter to concurrently track the multiple possible locations.

16
17 **30.** A method as recited in claim 27, further comprising using an
18 unscented particle filter to concurrently track the multiple possible locations.

19
20 **31.** A system to track multiple individuals in video content, the system
21 comprising:

22 an auto-initialization module to detect a candidate region for a new face in a
23 frame of the video content;

24 a hierarchical verification module to generate a confidence level for the
25 candidate region; and

1 a multi-cue tracking module to use a plurality of visual cues to track
2 previous candidate regions with confidence levels, generated by the hierarchical
3 verification module, that exceeded a threshold value.
4

5 **32.** A system as recited in claim 31, wherein the hierarchical
6 verification module is further configured to:

7 check whether the confidence level exceeds the threshold value;
8 if the confidence level does exceed the threshold value then to pass the
9 candidate region to the multi-cue tracking module; and

10 if the confidence level does not exceed the threshold value then to discard
11 the candidate region and not pass the candidate region to the multi-cue tracking
12 module.
13

14 **33.** A system as recited in claim 31, wherein the hierarchical
15 verification module is further configured to:

16 receive, from the multi-cue tracking module, an indication of a region;
17 verify whether the region is a face; and
18 return the region to the multi-cue tracking module for continued tracking
19 only if the region is verified as a face.
20

21 **34.** A system as recited in claim 31, wherein the system comprises a
22 video conferencing system.
23
24
25

1 35. A system as recited in claim 31, wherein the auto-initialization
2 module is further to:

3 detect whether there is motion in the frame;

4 if there is motion in the frame, then perform motion-based initialization to
5 identify the candidate region;

6 detect whether there is audio in the frame;

7 if there is audio in the frame, then perform audio-based initialization to
8 identify the candidate region; and

9 if there is neither motion in the frame nor audio in the frame, then use a fast
10 face detector to identify the candidate region.

11
12 36. A system as recited in claim 31, wherein the hierarchical
13 verification module is to use one or more hierarchical verification levels that
14 include a coarse level and a fine level, wherein the coarse level can verify whether
15 the new face is in the candidate area faster but with less accuracy than the fine
16 level.

17
18 37. One or more computer readable media having stored thereon a
19 plurality of instructions that, when executed by one or more processors, causes the
20 one or more processors to:

21 receive an indication of an area of a frame of video content;

22 use a first verification process to determine whether a human head is in the
23 area; and
24
25

1 if the first verification process verifies that the human head is in the area,
2 then indicate the area includes a face, and otherwise use a second verification
3 process to determine whether the human head is in the area.

4
5 **38.** One or more computer readable media as recited in claim 37,
6 wherein the first verification process and the second verification process
7 correspond to a plurality of hierarchical verification levels.

8
9 **39.** One or more computer readable media as recited in claim 38,
10 wherein the plurality of hierarchical verification levels comprise more than two
11 hierarchical verification levels.

12
13 **40.** One or more computer readable media as recited in claim 37,
14 wherein the first verification process is a coarse level process and the second
15 verification process is a fine level process, and wherein the coarse level process
16 can verify whether the human head is in the candidate area faster but with less
17 accuracy than the fine level process.

18
19 **41.** One or more computer readable media as recited in claim 37,
20 wherein the plurality of instructions to use the first verification process comprises
21 instructions that cause the one or more processors to:

22 generate a color histogram of the area;

23 generate an estimated color histogram of the area based on previous frames
24 of the video content;

1 determine a similarity value between the color histogram and the estimated
2 color histogram; and

3 verify that the candidate area includes the human head if the similarity
4 value is greater than a threshold value.

5
6 **42.** One or more computer readable media as recited in claim 37,
7 wherein the plurality of instructions to receive the indication of the area of the
8 frame of video content comprises instructions that cause the one or more
9 processors to:

10 receive a candidate area for a new face region in the frame.

11
12 **43.** One or more computer readable media as recited in claim 37,
13 wherein the plurality of instructions to receive the indication of the area of the
14 frame of video content comprises instructions that cause the one or more
15 processors to:

16 receive an indication of an area to re-verify as including a face.

17
18 **44.** One or more computer readable media having stored thereon a
19 plurality of instructions to detect a candidate region for an untracked face in a
20 frame of content, wherein the plurality of instructions, when executed by one or
21 more processors, causes the one or more processors to:

22 detect whether there is motion in the frame;

23 if there is motion in the frame, then perform motion-based initialization to
24 identify the candidate region;

25 detect whether there is audio in the frame;

1 if there is audio in the frame, then perform audio-based initialization to
2 identify the candidate region; and

3 if there is neither motion in the frame nor audio in the frame, then use a fast
4 face detector to identify the candidate region.

5
6 **45.** One or more computer readable media as recited in claim 44,
7 wherein the plurality of instructions to perform motion-based initialization
8 comprises instructions that cause the one or more processors to:

9 determine whether there is motion at a plurality of pixels on a plurality of
10 lines across the frame;

11 generate a sum of frame differences for a plurality of segments of multiple
12 ones of the plurality of lines;

13 select, for each of the multiple lines, the segment having the largest sum;

14 identify a smoothest region of the selected segments;

15 check whether the smoothest region resembles a human upper body; and

16 extract, as the candidate area, the portion of the smoothest region that
17 resembles a human head.

18
19 **46.** One or more computer readable media as recited in claim 45,
20 wherein the instructions to determine whether there is motion comprise
21 instructions that cause the one or more processors to:

22 determine, for each of the plurality of pixels, whether a difference between
23 an intensity value of the pixel in the frame and an intensity value of a
24 corresponding pixel in one or more other frames exceeds a threshold value.

1 47. One or more computer readable media having stored thereon a
2 plurality of instructions to track faces from frame to frame of content, wherein the
3 plurality of instructions, when executed by one or more processors, causes the one
4 or more processors to:

5 predict, using a plurality of cues, where a contour of a face will be in a
6 frame;

7 encode a smoothness constraint that penalizes roughness;

8 apply the smoothness constraint to a plurality of possible contour locations;

9 and

10 select the contour location having the smoothest contour as the location of
11 the face in the frame.

12
13 48. One or more computer readable media as recited in claim 47,
14 wherein the plurality of cues include foreground color, background color, edge
15 intensity, and motion.

16
17 49. One or more computer readable media as recited in claim 47,
18 wherein the plurality of cues include audio.

19
20 50. One or more computer readable media as recited in claim 47,
21 wherein the smoothness constraint includes contour smoothness.
22
23
24
25

1 **51.** One or more computer readable media as recited in claim 47,
2 wherein the smoothness constraint includes both contour smoothness and region
3 smoothness.

4
5 **52.** One or more computer readable media as recited in claim 47,
6 wherein the plurality of instructions to encode the smoothness constraint
7 comprises instructions that cause the one or more processors to generate Hidden
8 Markov Model (HMM) state transition probabilities.

9
10 **53.** One or more computer readable media as recited in claim 47,
11 wherein the plurality of instructions to encode the smoothness constraint
12 comprises instructions that cause the one or more processors to generate Joint
13 Probability Data Association Filter (JPDAF) state transition probabilities.

14
15 **54.** One or more computer readable media as recited in claim 47,
16 wherein the plurality of instructions further comprise instructions that cause the
17 one or more processors to:

18 adapt the predicting for the face in subsequent frames to account for
19 changing color distributions.

20
21 **55.** One or more computer readable media as recited in claim 47,
22 wherein the plurality of instructions further comprise instructions that cause the
23 one or more processors to:

24 adapt the predicting for the face in subsequent frames based on one or more
25 cues observed in the frame.

1
2 **56.** One or more computer readable media as recited in claim 47, the
3 plurality of instructions further comprise instructions that cause the one or more
4 processors to concurrently track multiple possible locations for the face from
5 frame to frame.

6
7 **57.** One or more computer readable media as recited in claim 56, the
8 plurality of instructions further comprise instructions that cause the one or more
9 processors to concurrently track the multiple possible locations.

10
11 **58.** A method for tracking an object along frames of content, the method
12 comprising:

13 using a plurality of cues to track the object.

14
15 **59.** A method as recited in claim 58, wherein the plurality of cues
16 include foreground color, background color, edge intensity, motion, and audio.

17
18 **60.** A method as recited in claim 58, wherein the using comprises
19 predicting wherein the object will be from frame to frame based on the plurality of
20 cues.

21
22 **61.** A method for tracking an object along frames of content, the method
23 comprising:

24 predicting where the object will be in a frame;

25 encoding a smoothness constraint that penalizes roughness;

1 applying the smoothness constraint to a plurality of possible object
2 locations; and

3 selecting the object location having the smoothest contour as the location of
4 the object in the frame.

5
6 **62.** A method as recited in claim 61, wherein the predicting uses a
7 plurality of cues that include foreground color, background color, edge intensity,
8 motion, and audio.

9
10 **63.** A method as recited in claim 61, wherein the smoothness constraint
11 includes both contour smoothness and region smoothness.

12
13 **64.** A method as recited in claim 61, wherein encoding the smoothness
14 constraint comprises generating Hidden Markov Model (HMM) state transition
15 probabilities.

16
17 **65.** A method as recited in claim 61, wherein encoding the smoothness
18 constraint comprises generating Joint Probability Data Association Filter (JPDAF)
19 state transition probabilities.

20
21 **66.** A method as recited in claim 61, wherein using the plurality of cues
22 to track each verified face further comprises, for each face:

23 adapting the predicting for the face in subsequent frames based on one or
24 more cues observed in the frame.

1 67. A method as recited in claim 61, wherein predicting where the
2 object will be comprises:

3 accessing a set of one or more feature points of the face; and

4 analyzing the frame to identify an area that includes the set of one or more
5 feature points.

6
7 68. A method as recited in claim 61, wherein using the plurality of cues
8 to track each verified face comprises concurrently tracking multiple possible
9 locations for the face from frame to frame.

10
11 69. A method as recited in claim 68, further comprising using a
12 multiple-hypothesis tracking technique to concurrently track the multiple possible
13 locations.

14
15 70. A method as recited in claim 61, wherein the object comprises a
16 face in video content.

17
18 71. A method as recited in claim 61, wherein the object comprises a
19 sound source location in audio content.